

RGB Images

IR Images

Encoder

Encoder

Embedding Spaces

Clusters

Matching

Similarity Matrix

Diverse Tokens Matching

 \mathcal{L}_{V2R}

V2R:

R2V:

 \mathcal{L}_{R2V}

Correspondences

Epoch $i+1$ Epoch i \mathcal{L}_{mate}

Homogeneous Fusion

Query

Local clustering by intra-camera training

 $\mathcal{N}(q_i)^{(1)}$ $\mathcal{N}(q_i)^{(2)}$ $\mathcal{N}(q_i)^{(4)}$ $L(q_i)$ $\mathcal{N}(q_i) \cap L(q_i) \cup \mathcal{N}(q_i)^{no-c}$

Optimization

 $\mathcal{L}_{Neighbor}$

Diverse Tokens Neighbor Learning

- ☐ / ☐ Different clusters with same ID
- ☐ / ☐ Visible/Infrared cluster
- ☐ Instance
- \longrightarrow Pull